

## Note

**Shortest consistent superstrings computable  
in polynomial time\***

Tao Jiang\*, Vadim G. Timkovsky

*Department of Computer Science and Systems, McMaster University, Hamilton, Ont., Canada L8S 4K1*

Received July 1994; revised September 1994

Communicated by M. Nivat

---

**Abstract**

The shortest consistent superstring problem is, given a set of positive strings and a set of negative strings, finding a shortest string including every positive string and no negative string as a substring. This problem is NP-hard and arises in DNA sequencing by hybridization. It is also an extension of the well-known shortest common superstring problem which corresponds to the case when the set of negative strings is empty. In this paper we show that a shortest consistent superstring can be found in polynomial time if (i) a longest common nonsuperstring for the set of negative strings exists or (ii) the number of positive strings is bounded and every symbol of the alphabet appears at the end of some negative string. In the case (i) a longest consistent superstring can also be found in polynomial time.

---

**1. Introduction**

Jiang and Li [5, 6] were perhaps the first to pose the *shortest consistent superstring* problem: Given a set  $P$  of *positive* strings and a set  $N$  of *negative* strings over an alphabet  $\Sigma$ , find a *shortest consistent (with  $N$ ) superstring* for  $P$ , i.e. a shortest string over  $\Sigma$  including every positive string and no negative string as a substring. Since  $N$  will remain invariable further, we will omit the expression in the brackets. Obviously, this problem has a solution only if  $P \cap N = \emptyset$ . Further we assume (restrictively) that  $P \cup N$  is *inclusion free*, i.e., no string in the set includes another as a substring. In particular, this union does not include the empty string. Besides, without loss of generality we assume that there is no one-symbol negative string, say,  $a$  because, due

---

\* Supported in part by NSERC Research Grant OGP0046613.

\* Corresponding author. Email addresses: {jiang, timko}@maccs.mcmaster.ca.

to the inclusion free requirement, no string in  $P \cup N$  except  $a$ , therefore, no string sought for, contains  $a$ . So,  $a$  can be deleted from  $\Sigma$  and  $N$ .

This problem arises in DNA sequencing by hybridization [1,11] and the PAC learning model [16] for DNA sequencing [6,10] and can be considered as a natural extension of the well-known *shortest common superstring* problem [2,9,14] which corresponds to the case when  $N = \emptyset$ . Since the latter is NP-hard, the former is also so and remains NP-hard even if  $N \neq \emptyset$  [7]. In the case  $P = \emptyset$ , the shortest consistent superstring problem becomes trivial. However, replacing the “shortest” requirement by the “longest” one transforms this case into the *longest common nonsuperstring* problem [15] solvable in polynomial time [12,13]. Consistent superstrings can be considered as intermediate structures between common superstrings and common nonsuperstrings. It prompts the conjecture that there are polynomially solvable nontrivial cases of the shortest consistent superstring problem that somehow are related to longest common nonsuperstrings. As it turns out, these cases exist indeed and arise from the graph representation of common nonsuperstrings that has been proposed in [13].

In this paper we show that a shortest consistent superstring can be found in polynomial time if (i) a longest common nonsuperstring for  $N$  exists or (ii)  $|P|$  is bounded and every symbol of  $\Sigma$  appears at the end of some string in  $N$ . As we will see further, in the case (i) a *longest consistent superstring* can also be found in polynomial time.

We hope that these results will find applications in creating more effective methods for DNA sequencing by hybridization. Besides, the case (i) provides an example of “downfall” phenomenon – an NP-hard problem has a polynomially solvable extension, which is of some theoretical interest.

## 2. Preliminaries

### 2.1. Graph representation of common nonsuperstrings

Here we cite some notation and results in [13]. Let  $\alpha$  be a string over an alphabet  $\Sigma$  and  $n$  be natural. Then  $\alpha^n$  denotes the concatenation of  $n$  copies of  $\alpha$ . The following condition is necessary for the existence of a longest common nonsuperstring for  $N$ .

$$\forall a \in \Sigma \exists n \geq 2: a^n \in N \quad (\text{a quadratical closure of } N).$$

Indeed, let this condition be false. Then either  $a \in N$  for some  $a \in \Sigma$  or  $a^n \notin N$  for all natural  $n$ . The former is not the case because of the original assumption, the latter implies that  $a^n$  is a common nonsuperstring for  $N$  and so there is no longest one.

Let  $|\alpha|$  be the length of  $\alpha$ ,  $\varepsilon$  be the empty string,  $|\varepsilon| = 0$ ,  $\text{Pref}_k \alpha$  and  $\text{Suff}_k \alpha$  be the prefix and the suffix of length  $k$  of  $\alpha$ , where  $k = 0, 1, \dots, |\alpha|$  and  $\text{Pref}_0 \alpha = \text{Suff}_0 \alpha = \varepsilon$ . We will write  $\text{Pref} \alpha$  and  $\text{Suff} \alpha$  if the length of prefixes and suffixes is immaterial.

The quadratical closure condition implies the following weaker one.

$$\forall a \in \Sigma \exists \alpha \in N: a = \text{Suff}_1 \alpha \quad (\text{a final closure of } N).$$

So, this condition is also necessary for the existence of a longest common nonsuperstring. In the following, we assume that  $N$  is finally closed.

For a set  $V$  of strings over  $\Sigma$  without  $\varepsilon$  define a directed graph  $G_V$  with vertex set  $V$  and arc set  $E$  determined by the rule:

$$(\alpha, \beta) \in E \Leftrightarrow \text{Suff}_{|\alpha|-1} \alpha = \text{Pref}_{|\beta|-1} \beta.$$

The arc  $(\alpha, \beta)$  arises when  $|\alpha| \leq |\beta| - 1$  and there is a string of length  $|\beta| + 1$  over  $\Sigma$  with prefix  $\alpha$  and suffix  $\beta$ . This string is denoted as  $[\alpha, \beta]$ . Note that  $\alpha, \beta$  is not necessarily in  $V$ . In particular, if  $|\alpha| = 1$ , then there are arcs from  $\alpha$  to all other vertices of  $G_V$  and  $\text{Suff}_{|\alpha|-1} \alpha = \beta$  implies  $(\alpha, \beta) \in E$ .

**Example 2.1.**  $G_{\text{English}}$  contains arcs (word, order), (a, part), (there, here), herein [word, order] = worder, [a, part] = apart, [there, here] = there.<sup>1</sup>

Let  $\alpha(i)$  be the symbol at the  $i$ th position of  $\alpha$ , i.e.  $\alpha = \alpha(1)\alpha(2) \dots \alpha(|\alpha|)$ , and  $V^\geq$  be the set of all common superstrings for  $V$ . A route in a graph is denoted by the sequence  $\mathcal{A} = (\alpha_1, \alpha_2, \dots, \alpha_{k-1}, \alpha_k)$  of its vertices such that for all  $i = 1, 2, \dots, k-1$   $(\alpha_i, \alpha_{i+1})$  is an arc of the graph.<sup>2</sup>  $|\mathcal{A}|$  will denote the length of  $\mathcal{A}$ , i.e.  $|\mathcal{A}| = k-1$ . For every route  $\mathcal{A} = (\alpha_1, \alpha_2, \dots, \alpha_{k-1}, \alpha_k)$  in  $G_V$  define the string

$$f(\mathcal{A}) = \alpha_1(1)\alpha_2(1) \dots \alpha_{k-1}(1)\alpha_k.$$

Obviously,  $f(\mathcal{A}) \in \{\alpha_1, \alpha_2, \dots, \alpha_{k-1}, \alpha_k\}N^\geq$ .

**Remark 2.1.** Let  $\Sigma^n$  be the set of all strings of length  $n$  over  $\Sigma$ . Then  $G_{\Sigma^n}$  is the well-known de Bruijn's graph [3] and  $f$  is a one-to-one correspondence between the set of routes of length  $k-1$  in  $G_{\Sigma^n}$  and  $\Sigma^{n+k-1}$ .

For a string  $\alpha$  with  $|\alpha| > 1$  we call  $\text{Pref}_n \alpha$  and  $\text{Suff}_n \alpha$  proper if  $0 < n < |\alpha|$ . Let  $S$  be the set of proper suffixes of strings in  $N$ . For every nonempty string  $\omega$  over  $\Sigma$  define the route in  $G_S$

$$g(\omega) = (\alpha_1, \alpha_2, \dots, \alpha_i, \alpha_{i+1}, \dots, \alpha_{|\omega|-1}, \alpha_{|\omega|}),$$

where  $\alpha_i$  is the longest string of  $S$  included in  $\omega$  as a substring starting from the  $i$ th position. Since the final closure condition holds,  $a \in \Sigma \Rightarrow a \in S$  and so the choice of  $\alpha_i$  is always possible. It is important to observe that the inequality  $|\alpha_i| \leq |\alpha_{i+1}| + 1$  and

<sup>1</sup>Note that *English* is actually not inclusion free. The inclusion freeness assumption is only needed in the next section.

<sup>2</sup>Unlike a path, a route can intersect itself.

the equality  $\text{Suff}_{|x_i|-1} x_i = \text{Pref}_{|x_{i+1}|-1} x_{i+1}$  follow from the fact that  $x_i$  and  $x_{i+1}$  are included in  $\omega$  as substrings starting from the  $i$ th and  $(i+1)$ th positions, respectively, and the longest length requirement. Thus, the arc  $(x_i, x_{i+1})$  exists in fact, i.e. the definition of the route  $g(\omega)$  is correct. It is easy to see that  $f(g(\omega)) = \omega$ .

Let  $\Gamma_S$  be a subgraph of  $G_S$  with vertex set  $S$  and arcs  $(\alpha, \beta)$  so that  $\alpha$  is the longest prefix of  $|x, \beta|$  contained in  $N \cup S$ , i.e. among suffixes of  $N$  there are no prefixes of  $|x, \beta|$  longer than  $\alpha$ . Note that this definition is independent of the final closure condition.

**Example 2.2.** If  $N = \text{English}$  then the arc  $(ove, venir)$  of  $G_S$  is not in  $\Gamma_S$ , since  $ove = \text{Suff}_3 \text{love}$ ,  $venir = \text{Suff}_5 \text{souvenir}$ ,  $[ove, venir] = ovenir$ , but  $oven = (\text{Pref}_4 ovenir) \in N$ . Neither is the arc  $(e, t)$  of  $G_S$  in  $\Gamma_S$  since  $e = \text{Suff}_1 \text{love}$ ,  $t = \text{Suff}_1 \text{let}$ ,  $[e, t] = et = \text{Pref}_2 et = (\text{Suff}_2 \text{net}) \in S$ . However the arc  $(ee, ea)$  is in  $\Gamma_S$  since  $ee = \text{Suff}_2 \text{tree}$ ,  $ea = \text{Suff}_2 \text{tea}$ ,  $[ee, ea] = eea$ , where  $eea \notin N \cup S \ni ee$ .

Let  $V^\#$  be the set of all common nonsuperstrings for  $V$ . The following lemma [13] shows that  $\Gamma_S$  is constructed so that there is a correspondence between the set of routes in  $\Gamma_S$  and  $N^\#$ .

**Lemma 2.1.** (a) If  $\mathcal{A}$  is a route in  $\Gamma_S$ , then  $f(\mathcal{A}) \in N^\#$ ; (b) if  $\omega \in N^\#$ , then  $g(\omega)$  is a route in  $\Gamma_S$ .

**Remark 2.2.** Note that the last vertex of the route  $g(\omega)$  is a one-symbol suffix. Besides, an arbitrary route  $\mathcal{A} = (x_1, \dots, x_k)$  in  $\Gamma_S$  can be extended to the route

$$\mathcal{A}' = (x_1, \dots, x_k, \text{Suff}_{|x_k|-1} x_k, \text{Suff}_{|x_k|-2} x_k, \dots, \text{Suff}_1 x_k).$$

Lemma 2.1 reduces a consideration of common nonsuperstrings for  $N$  to the analysis of the graph  $\Gamma_S$ . Bounded length of common nonsuperstrings for  $N$  means bounded length of routes in  $\Gamma_S$ , i.e.  $\Gamma_S$  is acyclic. Hence, we have the following theorem [13].

**Theorem 2.1.** If the graph  $\Gamma_S$  is acyclic, then  $\mathcal{A}$  is the longest path in it if and only if  $f(\mathcal{A})$  is a longest common nonsuperstring for  $N$ . If  $\mathcal{A} = (x_1, x_2, \dots, x_k)$  is a closed route in  $\Gamma_S$ , i.e.  $x_1 = x_k$ , then  $[\text{Pref}_{k-1} f(\mathcal{A})]^n \in N^\#$  for any natural  $n$  so there is no longest common nonsuperstring for  $N$ .

Thus, a necessary and sufficient condition for the existence of a longest common nonsuperstring for  $N$  is the acyclicity of  $\Gamma_S$ .

## 2.2. Wreaths

Let  $\mathcal{A} = (x_1, x_2, \dots, x_k)$  and  $\mathcal{B} = (\beta_1, \beta_2, \dots, \beta_l)$  be two routes and  $k \leq l$ . If  $x_k = \beta_1$ , then  $\mathcal{A} \circ \mathcal{B}$  will denote the composition of  $\mathcal{A}$  and  $\mathcal{B}$ , i.e.  $\mathcal{A} \circ \mathcal{B} =$

$(\alpha_1, \alpha_2, \dots, \alpha_k, \beta_2, \dots, \beta_l)$ . We call  $\mathcal{A}$  a *head* of  $\mathcal{B}$ , if  $\alpha_1 = \beta_1, \alpha_2 = \beta_2, \dots, \alpha_k = \beta_k$ , a *tail* of  $\mathcal{B}$ , if  $\alpha_1 = \beta_{l-k+1}, \alpha_2 = \beta_{l-k+2}, \dots, \alpha_k = \beta_l$ , a *subroute* of  $\mathcal{B}$ , if  $\alpha_1 = \beta_i, \alpha_2 = \beta_{i+1}, \dots, \alpha_k = \beta_{i+k}$  for  $1 \leq i \leq l-k$ .

We call  $\alpha_1$  and  $\alpha_k$  a *source* and a *terminal* of  $\mathcal{A}$ , respectively. If  $\alpha_1 = \beta_1, \alpha_2 = \beta_2, \dots, \alpha_{l-1} = \beta_{l-1}, \alpha_l \neq \beta_l$ , then  $\alpha_{l-1}$  is a *fork* of  $\mathcal{A}$  and  $\mathcal{B}$  with *prongs*  $\alpha_l$  and  $\beta_l$ . A set of routes we call a *bundle* if all routes in it have a common source – the *knot* of a bundle. A *terminal*, a *fork* and a *prong* of a bundle are the terminal of a route, the fork and the prong of a pair of routes in it, respectively.

A *wreath* is a bundle in which at least one prong of every fork is a terminal and no two routes have a common terminal. Obviously, a longest route  $\mathcal{L}$  in a wreath contains all its forks. A *core* of a wreath is the one-vertex route including the knot, if  $|\mathcal{L}| = 0$ , or the head of  $\mathcal{L}$  of length  $|\mathcal{L}| - 1$  if  $|\mathcal{L}| > 0$ . It is easy to see that a wreath has only one core and each of the wreath's terminals is connected with the core by one arc. Thus, a wreath can be determined by its core and the arcs ending in terminals.

### 3. Representation of consistent superstrings by routes in $\Gamma_S$

#### 3.1. Representation of positive strings

An obvious necessary condition for the existence of a consistent superstring for  $P$  is  $P \subseteq N^\#$ . Further we consider that it holds.

The results from the previous section show that any common nonsuperstring for  $N$ , in particular, any string in  $P$ , any substring of a string in  $P$  and any consistent superstring for  $P$  can be represented as a route of  $\Gamma_S$  if the final closure condition holds. However, such a route is not unique in general (cf. Remark 2.2).

Let  $\mathcal{A} = (\alpha_1, \alpha_2, \dots, \alpha_{k-2}, \alpha_{k-1}, \alpha_k)$  be a route in  $\Gamma_S$ ,  $\omega \in N^\#$ . We say that  $\mathcal{A}$  is a *representative* of  $\omega$  if  $\omega$  is a prefix of  $f(\mathcal{A}) = \alpha_1(1)\alpha_2(1)\dots\alpha_{k-2}(1)\alpha_{k-1}(1)\alpha_k$  and not a prefix of  $f'(\mathcal{A}) = \alpha_1(1)\alpha_2(1)\dots\alpha_{k-2}(1)\alpha_{k-1}$ . Obviously,  $\text{Suff}_{|\omega|-k+1}\omega = \text{Pref}_{|\omega|-k+1}[\alpha_{k-1}, \alpha_k]$ . Here we consider that  $\alpha_0 = \varepsilon$  and  $[\varepsilon, x] = x$  for any vertex  $x$  of  $\Gamma_S$ . Since  $\omega$  is longer than  $f'(\mathcal{A})$ ,  $\alpha_{k-1}$  is a proper prefix of  $\text{Suff}_{|\omega|-k+1}\omega$ , i.e. this suffix is longer than  $\alpha_{k-1}$ . Therefore,  $\omega \in S$  implies that the arc  $(\alpha_{k-1}, \alpha_k)$  is not in  $\Gamma_S$ . This in turn implies  $k = 1$ . This proves

**Lemma 3.1.** *Every proper suffix of any negative string, i.e. every vertex of  $\Gamma_S$ , has representatives only of length 0.*

Let  $W_\omega$  denotes the set of all representatives of  $\omega$ . This set is not empty because it includes at least a head of  $g(\omega)$  (cf. Remark 2.2).

**Theorem 3.1.** *If  $\omega \in P$ , then  $W_\omega$  is a wreath of  $\Gamma_S$ .*

**Proof.** Let  $\mathcal{A} = (\alpha_1, \dots, \alpha_k), (\beta_1, \dots, \beta_l) \in W_\omega$ . Without loss of generality we assume that  $\alpha_1 = \text{Pref} \beta_1$ , since both  $\alpha_1$  and  $\beta_1$  are prefixes of  $\omega$ . If  $f(\mathcal{A}) = \text{Pref} \beta_1$ , then  $\omega \approx \text{Pref} \beta_1$ , i.e.  $P \cup N$  is not inclusion free – a contradiction with the original assumption. If  $\beta_1 = \text{Pref} f(\mathcal{A})$ , then either  $\alpha_1 = \beta_1$  or there exists the representative  $(\alpha_1, \dots, \alpha_i) \in W_{\beta_1}$  with  $i > 1$  – a contradiction with Lemma 3.1. Thus, all representatives of  $\omega$  have a common source, i.e.  $W_\omega$  is a bundle.

Now we show that at least one prong of each fork in  $W_\omega$  is a terminal. Let, in contrast,  $\phi_1$  and  $\psi$  be two prongs of a fork that are not terminals. Without loss of generality we assume that  $\phi_1 = \text{Pref} \psi$ . Then there exists the representative  $(\phi_1, \dots, \phi_i) \in W_\psi$  with  $i > 1$  – again a contradiction with Lemma 3.1.

Employing the same technique we can make sure that no two routes in  $W_\omega$  have a common terminal. Otherwise this terminal should have a representative of length more than 0.  $\square$

Further, if  $\omega \in P$ , we denote the knot and the core of the wreath  $W_\omega$  as  $\kappa_\omega$  and  $\mathcal{C}_\omega$ , respectively. Note that  $g(\omega)$  is the longest route in  $W_\omega$  and  $|g(\omega)| = |\omega| - 1$ . So,  $|\mathcal{C}_\omega| = 0$ , if  $|\omega| = 1$ , or  $|\mathcal{C}_\omega| = |\omega| - 2$  if  $|\omega| > 1$ .

### 3.2. Consistent superstrings for a pair of positive strings

Let  $\alpha, \beta \in P$ . If  $\text{Suff}_k \alpha = \text{Pref}_k \beta$ , where  $0 < k < \min\{|\alpha|, |\beta|\}$ , we say that the string  $(\text{Pref}_{|\alpha|-k} \alpha)\beta = \alpha \text{Suff}_{|\beta|-k} \beta$  is the *overlap* of  $\alpha$  and  $\beta$  with size  $k$ .

Since  $P$  is inclusion free, some consistent superstrings for the pair  $\{\alpha, \beta\}$  includes  $\alpha$  starting before  $\beta$  and the others have  $\beta$  before  $\alpha$ . So, we can talk just about consistent superstrings for the ordered pair  $\langle \alpha, \beta \rangle$  because all reasoning for  $\langle \beta, \alpha \rangle$  will be symmetrical.

We call a consistent superstring  $\gamma$  for  $\langle \alpha, \beta \rangle$  *proper* if  $\alpha = \text{Pref} \gamma$  and  $\beta = \text{Suff} \gamma$ . Obviously, any consistent superstring includes some proper one as a substring. Theorem 3.1 has the following evident.

**Corollary 3.1.** Any representative of a proper consistent superstring for  $\langle \alpha, \beta \rangle$  includes a route of  $W_\alpha$  as a head and a route of  $W_\beta$  as a tail.

So, any representative of a proper consistent superstring  $\gamma$  for  $\langle \alpha, \beta \rangle$  can be written as  $\mathcal{A} \circ \mathcal{B}$ , where  $\mathcal{A}$  is a route from  $\kappa_\alpha$  to  $\kappa_\beta$  and  $\mathcal{B}$  is a route of  $W_\beta$ . If  $\mathcal{A}$  includes a terminal  $\tau$  of  $W_\alpha$ , then  $\mathcal{A} = \mathcal{X} \circ \mathcal{Y}$ , where  $\mathcal{X}$  is a route of  $W_\alpha$  with the terminal  $\tau$  and  $\mathcal{Y}$  is a route from  $\tau$  to  $\kappa_\beta$ . If  $\mathcal{A}$  includes no terminal of  $W_\alpha$ , then  $\kappa_\beta$  is on the core of  $W_\alpha$ . Therefore, a head  $\mathcal{T}$  of  $\mathcal{B}$  is a tail of a route of  $W_\alpha$ , i.e.  $\gamma$  is an overlap of  $\alpha$  and  $\beta$  with size  $|\mathcal{T}|$ . Hence, every route of  $W_\beta$  contains a tail of a route of  $W_\alpha$  as a head. Thus, we have proved

**Theorem 3.2.** Let  $S_{\alpha\beta}$  be the set of all routes in  $\Gamma_S$  from  $\kappa_\alpha$  to  $\kappa_\beta$  that have a route of  $W_\alpha$  as a head,  $\mathcal{T}_{\alpha\beta}$  be the head of  $\mathcal{C}_\alpha$  from  $\kappa_\alpha$  to  $\kappa_\beta$  (if  $\kappa_\beta$  is on  $\mathcal{C}_\alpha$ ) and  $R_{\alpha\beta} = S_{\alpha\beta} \cup \{\mathcal{T}_{\alpha\beta}\}$ . Then

$\{\gamma_{\mathcal{A}}: \mathcal{A} \in R_{\alpha\beta}\}$ , where  $\gamma_{\mathcal{A}} = [\text{Pref}_{|\mathcal{A}|}f(\mathcal{A})]\beta$ , is the set of all proper consistent superstrings for  $\langle \alpha, \beta \rangle$ . Herein  $W_{\alpha\beta} = \{\mathcal{A} \circ \mathcal{B}: \mathcal{B} \in W_{\beta}\}$  is the wreath of all representatives of  $\gamma_{\mathcal{A}}$ .

Denote by  $S'_{\alpha\beta}$  the subset of  $S_{\alpha\beta}$  that consists of routes whose tails starting at terminals of routes in  $W_{\alpha}$  are actually paths, and  $R'_{\alpha\beta} = S'_{\alpha\beta} \cup \{\mathcal{T}_{\alpha\beta}\}$ . Note that the length of each route in  $R'_{\alpha\beta}$  is bounded. Define a distance from  $\alpha$  to  $\beta$  by the following two ways. Let

$$\text{Min}(\alpha, \beta) = \min \{|\text{Pref}_{|\mathcal{A}|}f(\mathcal{A})|: \mathcal{A} \in R'_{\alpha\beta}\}$$

and  $\text{Max}(\alpha, \beta)$  be the same expression after replacing “min” by “max”. Obviously,  $|\gamma_{\mathcal{A}}| = |\text{Pref}_{|\mathcal{A}|}f(\mathcal{A})| + |\beta|$ . Since a shortest consistent superstring is proper, Theorem 3.2. has the following

**Corollary 3.2.** *If  $\mathcal{A}$  is a shortest route in  $R'_{\alpha\beta}$ , then  $\gamma_{\mathcal{A}}$  is a shortest consistent superstring for  $\langle \alpha, \beta \rangle$  with length  $\text{Min}(\alpha, \beta) + |\beta|$ .*

Let  $\mathcal{M}$  and  $\mathcal{N}$  be longest paths in  $\Gamma_S$  to  $\kappa_{\alpha}$  and from  $\kappa_{\beta}$ , respectively. Since  $\mathcal{N}$  finishes with a one-symbol string (cf. Remark 2.2) and any consistent superstring includes a proper one as a substring, we have

**Corollary 3.3.** *If  $\Gamma_S$  is an acyclic graph and  $\mathcal{A}$  is a longest path in  $R'_{\alpha\beta}$ , then*

$$[\text{Pref}_{|\mathcal{M}|}f(\mathcal{M})][\text{Pref}_{|\mathcal{A}|}f(\mathcal{A})]f(\mathcal{N})$$

*is a longest consistent superstring for  $\langle \alpha, \beta \rangle$  with length  $|\mathcal{M}| + \text{Max}(\alpha, \beta) + |\mathcal{N}| + 1$ .*

### 3.3. Consistent superstrings for a linearly ordered set of positive strings

Let  $P = \{\alpha_1, \dots, \alpha_p\}$ . Since  $P$  is inclusion free, any consistent superstring for  $P$  includes positive strings in some linear order according to the appearances of their first symbols. So, we can say about consistent superstrings for a linearly ordered set of positive strings. Besides, the set of all consistent superstrings for  $P$  is the union of the sets of consistent superstrings for all  $p!$  linear order of  $P$ . In this subsection we consider consistent superstrings for a fixed linear order  $\langle \alpha_1, \dots, \alpha_p \rangle$ . An evident corollary from Theorem 3.1 is

**Corollary 3.4.** *If a consistent superstring  $\gamma$  exists for  $\langle \alpha_1, \dots, \alpha_p \rangle$ , then any representative of  $\gamma$  includes a route of  $W_{\alpha_i}$ ,  $i = 1, \dots, p$ , as a subroute in the  $i$ th order.*

Analogously, we call a consistent superstring  $\gamma$  for  $\langle \alpha_1, \dots, \alpha_p \rangle$  proper if  $\alpha_1 = \text{Pref}\gamma$  and  $\alpha_p = \text{Suff}\gamma$ . So, any representative of a proper consistent superstring  $\gamma$  for  $\langle \alpha_1, \dots, \alpha_p \rangle$  can be written as  $\mathcal{A} \circ \mathcal{A}_2 \circ \dots \circ \mathcal{A}_{p-1} \circ \mathcal{B}$ , where  $\mathcal{A}_i \in R_{\alpha_{i+1}}$ ,  $1 \leq i \leq p$ , and  $\mathcal{B}$  is a route of  $W_{\alpha_p}$ . Thus, we have

**Theorem 3.3.** Let  $\delta_{\mathcal{A}_i} = \text{Pref}_{[\mathcal{A}_i]} f(\mathcal{A}_i)$  for  $i = 1, 2, \dots, p-1$  and  $\gamma_{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_{p-1}} = \delta_{\mathcal{A}_1} \delta_{\mathcal{A}_2} \dots \delta_{\mathcal{A}_{p-1}} x_p$ . Then

$$\{\gamma_{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_{p-1}} : \mathcal{A}_1 \in R_{x_1 x_2} \text{ \& } \mathcal{A}_2 \in R_{x_2 x_3} \text{ \& } \dots \text{ \& } \mathcal{A}_{p-1} \in R_{x_{p-1} x_p}\}$$

is the set of all proper consistent superstrings for  $\langle x_1, \dots, x_p \rangle$ . Herein

$$W_{\gamma_{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_{p-1}}} \{ \mathcal{A}_1 \circ \mathcal{A}_2 \circ \dots \circ \mathcal{A}_{p-1} \circ \mathcal{B} : \mathcal{B} \in W_{x_p} \}$$

is the wreath of all representatives of  $\gamma_{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_{p-1}}$ .

Analogously, denoting by  $\mathcal{M}$  and  $\mathcal{N}$  longest paths in  $\Gamma_S$  to  $\kappa_{x_1}$  and from  $\kappa_{x_p}$ , respectively, we get the following corollaries from Theorem 3.3.

**Corollary 3.5.** If  $\mathcal{A}_i$  is a shortest route in  $R'_{x_i x_{i+1}}$ ,  $1 \leq i < p$ , then  $\gamma_{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_{p-1}}$  is a shortest consistent superstring for  $\langle x_1, \dots, x_p \rangle$  with length

$$\sum_{i=1}^{p-1} \text{Min}(x_i, x_{i+1}) + |x_p|.$$

**Corollary 3.6.** If  $\Gamma_S$  is an acyclic graph and  $\mathcal{A}_i$  is a longest path in  $R'_{x_i x_{i+1}}$ ,  $1 \leq i < p$ , then  $[\text{Pref}_{[\mathcal{M}]} f(\mathcal{M})] \delta_{\mathcal{A}_1} \delta_{\mathcal{A}_2} \dots \delta_{\mathcal{A}_{p-1}} f(\mathcal{N})$  is a longest consistent superstring for  $\langle x_1, \dots, x_p \rangle$  with length

$$|\mathcal{M}| + \sum_{i=1}^{p-1} \text{Max}(x_i, x_{i+1}) + |\mathcal{N}| + 1.$$

#### 4. Polynomially computable cases

##### 4.1. The case when the number of positive strings is bounded and the final closure condition holds

Let, as before,  $p = |P|$  and  $n$  be the total length of negative strings. For any  $\omega \in P$ , the core  $\mathcal{C}_\omega$  can be obtained from the route  $g(\omega)$  by deleting the terminal. Therefore,  $W_\omega$  can be extracted from  $\Gamma_S$  by finding the route  $g(\omega)$  and testing if every vertex of  $\Gamma_S$  is a terminal connected with  $\mathcal{C}_\omega$  by one arc. A route in  $S'_{\alpha\beta}$  can be written as  $\mathcal{A} \circ \mathcal{P}$ , where  $\mathcal{A} \in W_\alpha$  and  $\mathcal{P}$  is a path from the terminal of  $\mathcal{A}$  to the knot  $\kappa_\beta$ . So, finding a shortest or longest route in  $R'_{\alpha\beta}$  reduces to enumerating terminals of  $W_\alpha$  and finding a shortest or longest path. Thus, if the final closure condition holds, then checking the existence of a consistent superstring and finding a shortest one for a linearly ordered set of positive strings reduce to constructing the graph  $\Gamma_S$ , extracting the wreaths and, by Corollary 3.5, finding at most  $p-1$  shortest paths. All these manipulations take time at most  $O(pn^2)$ . Hence, a shortest consistent superstring for  $P$  can be found by enumerating all  $p!$  linear orders on  $P$  in time  $O(p!pn^2)$  which is polynomial if  $p$  is bounded by a constant.



#### 4.2. The case when a longest common nonsuperstring for the set of negative strings exists

In this case the final closure condition holds and, by Theorem 2.1,  $\Gamma_S$  is an acyclic graph. Introduce a binary relation  $<$  on  $P$  putting  $\alpha < \beta$  if and only if there is a path in  $\Gamma_S$  from  $\kappa_\alpha$  to  $\kappa_\beta$ . Obviously,  $<$  is transitive and, since  $\Gamma_S$  is an acyclic graph, antisymmetrical. If a consistent superstring  $\gamma$  exists for  $P$ , then the inclusion free requirement for  $P$  implies that no two positive strings start in  $\gamma$  with the same position. Hence,  $\kappa_\alpha = \kappa_\beta \Rightarrow \alpha = \beta$ , i.e.  $<$  is a strict order. Besides, for every pair of positive strings  $\alpha$  and  $\beta$  the knots  $\kappa_\alpha$  and  $\kappa_\beta$  are connected by a subroute of  $g(\cdot)$ , i.e. every two positive strings are comparable by  $<$ . Hence, we have

**Theorem 4.1.** *Let a longest common nonsuperstring for  $N$  exist. If a consistent superstring for  $P$  exists, then  $<$  is a linear order and the set of all consistent superstrings for  $P$  is the set of all consistent superstrings for the linearly ordered set  $P$  with  $<$ .*

Identifying a linear order on  $P$  takes time at most  $O(p^2)$  [8]. Thus, using Corollaries 3.5 and 3.6, we can find a shortest and longest consistent superstring for  $P$  in polynomial time  $O(p + pn^2)$ .

### 5. Concluding remarks

In the case when the final closure condition holds, the shortest consistent superstring problem can be easily reduced to the *minimum weighted directed Hamiltonian path* problem: Find a path with minimum weight that goes through all vertices of a given complete directed graph with weighted arcs. Without loss of generality, it is enough to consider that the first and last vertices of the path are also given.

Let  $H$  be the complete directed graph with vertex set  $\{\kappa_\omega: \omega \in P\} \cup \{s, t\}$ . Define weights  $\omega$  on arcs  $(\kappa_\alpha, \kappa_\beta)$  of  $H$  by the following way. Put

$$\omega(\kappa_\alpha, \kappa_\beta) = \begin{cases} \text{Min}(\alpha, \beta) & \text{if } R'_{\alpha, \beta} \neq \emptyset, \\ \infty & \text{otherwise,} \end{cases}$$

and  $\omega(s, \kappa_\omega) = 0$ ,  $\omega(\kappa_\omega, t) = |\omega|$ . Corollary 3.5 implies that, if  $\mathcal{H} = (s, x_1, \dots, x_p, t)$  is a minimum weighted directed Hamiltonian path in  $H$  from  $s$  to  $t$ , then  $\gamma_{x_1, \dots, x_{p-1}}$  is a shortest consistent superstring. A dynamic programming algorithm computes  $\mathcal{H}$  in time  $\Theta(p^2 2^p)$  (cf. [4]). So, if  $p = O(\log n)$ , then a shortest consistent superstring can be found in polynomial time  $O(n^{O(1)} \log^2 n + n^2 \log n)$ . Besides, any efficiently solvable case of the minimum weighted directed Hamiltonian path problem gives an efficiently solvable case of the shortest consistent superstring problem if the final closure condition holds. In conclusion, one has to note that a key construction in all reasonings in this paper is the route  $g(\omega)$  which does not exist in general if the final closure condition does not hold.

## References

- [1] R. Drmanac and C. Crkvenjakov, Sequencing by hybridization (SBH) with oligonucleotide probes as an integral approach for the analysis of complex genomes, *Internat. J. Genomic Res.* 1 (1) (1992) 59–79.
- [2] J. Gallant, D. Maier and J. Storer, On finding minimal length superstrings, *J. Comput. System Sci.* 20 (1980) 50–58.
- [3] M. Hall Jr., *Combinatorial Theory* (Blaisdell, Waltham, MA, 1967).
- [4] E. Horowitz and S. Sahni, *Fundamentals of Computer Algorithms* (Computer Science Press, Rockville, MD, 1978).
- [5] T. Jiang and M. Li, Approximating shortest superstrings with constraints, *Theoret. Comput. Sci.* 134 (1994) 473–491.
- [6] T. Jiang and M. Li, DNA sequencing and string learning, *Math. Systems Theory*, to appear.
- [7] T. Jiang and M. Li, On the complexity of learning strings and sequences, *Theoret. Comput. Sci.* 119 (1993) 363–371.
- [8] D.E. Knuth, *The Art of Computer Programming, Vol. 3: Sorting and Searching* (Addison-Wesley, Reading, MA, 1973).
- [9] A. Lesk, ed., *Computational Molecular Biology, Sources and Methods for Sequence Analysis* (Oxford University Press, Oxford, 1988).
- [10] M. Li, Towards a DNA sequencing theory, in: *Proc. 31st IEEE Symp. on Foundations of Computer Science* (1990) 125–134.
- [11] P. Pevzner and R. Lipshutz, Towards DNA sequencing by hybridization, manuscript, 1993.
- [12] A.R. Rubinov and V.G. Timkovsky, String non-inclusion optimization problems, *SIAM J. Discrete Math.*, submitted.
- [13] A.R. Rubinov and V.G. Timkovsky, Non-similarity combinatorial problems I, *Biosystems* 30 (1993) 81–92; II, in: P.A. Pevzner and M.S. Gelfand, eds., *Computer Genetics* (Elsevier, Amsterdam, 1993) 81–92.
- [14] J. Storer, *Data Compression: Methods and Theory* (Computer Science Press, Rockville, MD, 1988).
- [15] V.G. Timkovskii, Complexity of common subsequence and supersequence problems and related problems, *Kibernetika* (1989) (5) 1–13; an English translation appears in: *Cybernetics* 25 (1990) 565–580.
- [16] L.G. Valiant, A theory of the learnable, *Comm. ACM* 27 (1984) 1134–1142.